



TITLE:

形式言語における区別・説明の無限過程の定式化 (理論計算機科学の深化と応用)

AUTHOR(S):

植村, 仁

CITATION:

植村, 仁. 形式言語における区別・説明の無限過程の定式化 (理論計算機科学の深化と応用). 数理解析研究所講究録 2009, 1649: 105-112

ISSUE DATE:

2009-05

URL:

<http://hdl.handle.net/2433/140749>

RIGHT:

形式言語における区別・説明の無限過程の定式化

真理大學資訊工程學系 植村 仁 (Jin Uemura)

Department of Computer Science and Information Engineering,
Aletheia University

1 導入

演繹, 帰納, 類推, 仮説演繹, 分類など, 思考のさまざまな面が定式化され, 研究されてきた. この研究では, 区別と説明に光を当て, これに複数の定式化を与え, それらがどのような性質をもつかを解明する.

入力例を説明する言語を出力してゆく過程については研究はなされてきたが, 例と言語を入力し, 事実を区別してゆく過程, つまり, 事実とルールを学んでゆき, 事実を区別し, かつルールを取捨選択してゆく過程については研究されていない. これを解明する.

2 基本的定義

2.1 概略と基本的定義

本稿における学習機械とは, 直感的には学習機械にとって未知の対象言語 F から出てくる語の例の無限列と, 入力例の説明と区別に用いる言語族 \mathcal{C} の言語の列 (正確には, 各言語を表す記述: オートマトンや文法などの列) を受け取り, 出力は入力される言語の列の部分からなり, 入力例を説明・区別するために十分なものを選ぶ. この選択基準については様々な定義が考えられる.

記号等の定義 アルファベット (定数記号の有限集合) を Σ とする. 本稿では, 言語はすべて Σ 上の言語であるとする.

入力される語についての定義等 学習対象言語の族を \mathcal{F} , その各言語を F などで表す. 言語 F の正提示とは, 語の無限列 e_0, e_1, e_2, \dots で $\{e_i \mid i \in \mathbb{N}\} = F$ を満たすものである. 主に σ などで表す. 正提示 $\sigma = e_0, e_1, e_2, \dots$ の n -初期断片とは, $e_0, e_1, e_2, \dots, e_n$ のことであり, $\sigma[n]$ と表す. F の部分で入力された既知の例を表す集合を E などで表す.

入力例を区別する言語などに関する定義等 入力例を区別するために用いる言語の族を \mathcal{C} , その言語族の各言語の記述 (オートマトンや文法など) がなす集合を $\mathcal{C}d$ などで表す. 本稿では, この区別に用いる言語族は帰納的言語の添え字つき族であると仮定する. つまりその言語族が L_0, L_1, L_2, \dots と表され, ある帰納的関数 $f: \mathbb{N} \times \Sigma^* \rightarrow \{0, 1\}$ が存在し, $f(i, w) = 1 \iff w \in L_i$ となることを仮定する. この自然数は言語の記述をゲーデル化したものとみなす. 言語族 \mathcal{C} , 言語族の記述の集合 $\mathcal{C}d$ の正提示, その n -初期断片は言語の正提示の場合にならって定義する. 入力された既知の言語の族 \mathcal{K} , その記述からなる集合を $\mathcal{K}d$ などで表す.

2.2 例で使用する言語族: 正規パターン言語族

変数の加算無限集合を X としよう. 正規パターンとは $(X \cup \Sigma)$ 上の有限文字列である. また, 正規パターン $p = w_0 x_1 w_1 \cdots x_n w_n$ ($w_0, w_n \in \Sigma^*$, $w_i \in \Sigma^+$ ($i = 1, \dots, n-1$), $x_i \in X$, 変数はすべて異なる) の言語 $L(p)$ とは, $\{w_0\}\Sigma^*\{w_1\}\cdots\Sigma^*\{w_n\}$ である. つまり $L(p)$ は接頭語 w_0 , 接尾語 w_n を共通して持ち, さらに重複なく w_1, \dots, w_{n-1} が出現するような語の集合であるともいえる. $w_0, w_n = \varepsilon$ (空語) であり, $w_1, \dots, w_{n-1} \in \Sigma$ であれば, $L(p)$ は w_1, \dots, w_{n-1} を部分列 (subsequence) としてもつような語の集合となる. このような正規パターンは部分列パターンとよぶことにしよう.

2.3 学習機械と例の区別

説明体系 本論文では既知の例を説明する際, 既知の言語の内どれを用いるかに制限を設ける. 既知の例と言語に対し, どのような言語を選ぶことが許されるかということ関数を用いて定義しよう. 形式的には,

定義 2.1 (説明体系). 説明体系とは $2^C \times 2^{\Sigma^*} \rightarrow 2^C$ を満たす関数であるとする. PT 等で表す.

次節において, 様々な説明体系を定義しその性質について議論するが, ここではその一つ目のものを例として取り上げよう.

例 2.2. \mathcal{K} を言語族とすると, 説明体系 PT_{dist} を $PT_{dist}(\mathcal{K}, E) = 2^{\mathcal{K}}$ としよう. この説明体系は, 既知の例の集合 E に関係なく, 既知の言語の族 \mathcal{K} の任意の部分許すようなものである.

学習機械 ここで定義する学習機械は, 未知の言語からの例と, 例を区別し説明するための言語の記述を入力として受け取り, 説明体系で許された入力言語の記述の部分を出力するというプロセスを無限に繰り返すものである.

定義 2.3 (学習機械). 説明体系 PT , 言語族 \mathcal{C} 上の学習機械 M とは, 入力と出力を限りなく続けるアルゴリズムである. 以下のように二種の初期断片を入力とし, PT が許す値をとる.

入力 1: 学習者にとって未知の対象言語族 \mathcal{F} の未知の言語を F とするとき, F の正提示の初期断片. ここでは $\sigma[n] = e_0, e_1, \dots, e_n$ と表しておく.

入力 2: 区別言語の族 \mathcal{C} とそれを表す記述の集合を Cd とするとき, Cd のある正提示. ここでは n' -初期断片 $\sigma'[n'] = d_0, d_1, \dots, d_{n'}$ と表すとする.

出力: $L(d_i)$ を d_i が表す言語, $\mathcal{K}_{n'} = \{L(d_0), L(d_1), \dots, L(d_{n'})\}$, $E_n = \{e_0, e_1, \dots, e_n\}$ とするとき, $T \in PT(\mathcal{K}_{n'}, E_n)$ となるある T である.

例の分割 出力の表す言語の族 T は, 既知の例の集合 E を分割することができる. T による E の分割をまず定義し, 次により細かい分割, 最後に極大な分割という概念を定義しよう.

定義 2.4 (分割). $T \subseteq \mathcal{C}$, $E \subseteq \Sigma^*$

$$D(T, E) = \{S \subseteq E \mid \begin{aligned} &S \neq \emptyset, \\ &S = E \cap (\bigcap_{L \in P} L) \cap (\bigcap_{L \in N} L^C), \\ &P \cup N = T, P \cap N \neq \emptyset \end{aligned}\}$$

分割の元を類と呼ぶことにしよう.

定義 2.5 (より細かい分割). $E \subseteq \Sigma^*$, $D, D' \subseteq 2^{\Sigma^*}$, D, D' はそれぞれ共通部分を持たない Σ^* の部分集合からなり, それぞれの和が E と一致とするとする.

D' が D のより細かい E の分割であるとは, 任意の $S \in D$ に対し, S が D' の元の和となることであるとする.

定義 2.6 (極大分割). PT を説明体系, \mathcal{K} を既知の言語族, E を既知の例集合とするとき, 分割 D (または $D = D(T, E)$ となる T) が (PT, \mathcal{K}, E) における極大分割であるとは, 任意の $T' \in PT(\mathcal{K}, E)$ に対し, $D(T', E)$ が D より細かい E の分割とならないことである.

例 2.7. 区別するための言語族の記述の集合 Cd として部分列パターンの集合をとり, 未知の言語 F の例 $E = \{aa, ac, abb, bbb\}$ と Cd の部分 $Kd = \{xay, xbbby, xby\}$ を既に受け取ったとする. \mathcal{K} を Kd の言語の族としておく.

説明体系 $PT_{dist}(\mathcal{K}, E) = 2^{\mathcal{K}} - \{\emptyset\}$ を用いた場合, $(PT_{dist}, \mathcal{K}, E)$ における極大分割となる $PT_{dist}(\mathcal{K}, E)$ の元は, $\{L(xay), L(xbbby), L(xby)\}$, $\{L(xay), L(xbbby)\}$, $\{L(xay), L(xby)\}$ であり, これらの一つを T とすると, $D(T, E) = \{\{aa, ac\}, \{abb\}, \{bbb\}\}$ と分割される.

Kd に xcy を追加し極大分割をすると, $PT_{dist}(\mathcal{K}, E)$ の元は, $\{L(xay), L(xbbby), L(xby), L(xcy)\}$, $\{L(xay), L(xbbby), L(xcy)\}$, $\{L(xay), L(xby), L(xcy)\}$ であり, 前段階の分割の類 $\{aa, ac\}$ が分割されることになる.

さらに E に b を追加すると, 極大分割するために今まで一つあればよかった xby と $xbby$ の双方が必要となる.

2.4 学習可能性の定義

区別と説明の無限過程に対し学習可能性の定義は複数ありうるが, ここでは極限同定モデル ([1]) に類似した学習可能性を定義する. 若干敷衍するならばその定義は, 学習機械が出力する言語族の記述の集合がある時点から変化せず, それが永久に例を極大に分割しつづけるということに基づいたものである. 本稿ではこの定義を用い議論を進める.

定義 2.8 (極限における極大分割). 学習機械の出力が収束するとは, 出力の無限列を Td_i ($i \in \mathbb{N}$) とするとき, ある $t \in \mathbb{N}$ 以降の出力が変化しないことである. つまり, 任意の $j \geq t$ に対し $Td_j = Td_t$ となることである.

説明体系 PT , 言語族 \mathcal{C} 上の学習機械が言語 F を極限において極大分割するとは, 出力が Td_t に収束し, かつ以下の条件を満たすことである: ステップ t における出力を Td_t , その言語族を T_t 既知の区別言語族の記述の集合を Kd_t , その言語族を \mathcal{K}_t 例の集合を E_t とするとき, $\forall \mathcal{K} (\mathcal{K}_t \subset \mathcal{K} \subseteq \mathcal{C})$, $\forall E (E_t \subset E \subseteq F)$, $\forall T \in PT(\mathcal{K}, E)$ に対し, $D(T_t, E)$ は (PT, \mathcal{K}, E) における極大分割となることである.

定義 2.9 (学習可能性). PT を説明体系とする, 言語族 \mathcal{C} の下で, 言語 F が PT 学習可能であるとは, 説明体系 PT , 言語族 \mathcal{C} 上のある学習機械 M が存在し, それが学習機械が F を極限において極大分割することである.

言語族 \mathcal{C} の下で, 言語族 \mathcal{F} が PT 学習可能であるとは, \mathcal{F} の任意の言語が説明体系 PT , 区別言語 \mathcal{C} の下で, 学習可能であることである.

言語族 \mathcal{C} の下で, \mathcal{C} が PT 学習可能であることを, ただ単に \mathcal{C} が PT 学習可能であるということにする.

3 4つの説明体系

学習機械は例を区別するための言語を次々と受け取る。受け取られた既知の言語の族の一部分で、入力された例を説明・区別する言語をいくつか選びだす。説明体系はその許される出力の定義となるが、さまざまな説明体系が考えられる。

3.1 言語・例・説明・内包・外延

まず学習機械の出力の表す言語族が説明する例の集合の族を定義し、それに注目し、いくつかの説明体系を定義してみる。

例の集合 E の一部分 S を説明する L があるとするとき、集合の内包的・外延的定義という用語に倣い、 S と L の関係を定義しよう。

定義 3.1 (内包・外延). L を言語, E, S を語の集合とする, $S = L \cap E$ であるとき, L を S の内包, S を L の外延と呼び, $S = Ex(L, E)$ と書くことにする。

一般には、つまり L, E が言語であるというだけであるならば、 L, E の積が計算可能であればよいが、本稿では、 E は有限、区別のための言語族の元 L は帰納的集合と仮定しているので、特に何も条件は必要ない。

定義 3.2 (外延の族). E を例となる語の集合, T を言語の族とする。外延の族を

$$Ex(T, E) = \{Ex(L, E) \mid L \in T\}$$

と書くことにする。

この外延の族を用いて、集合の交わりと包含関係を軸に、4つの説明体系を定義する。

1. 外延の族に特に条件のないもの。
2. 外延同士が、常に包含関係にあるもの。
3. 外延どうしが交わりを持たないもの。
4. 外延が包含関係にあるか、交わりを持たないようなもの。

3.2 完全区別

最初に挙げる説明体系は、入力される言語の記述の任意の部分を出し示るようなものである。直感的には、既知の言語で区別できるものはすべて区別しうる、というものである。この定義は、二階以上の述語論理における同一性 (equality) の定義 ($x = y \iff \forall P (P(x) \leftrightarrow P(y))$) にヒントを得たものである。この完全区別を用いた時、その極大区別の同一の類に属する二つの語は、既知の区別言語により分断されない、という性質を持つ。

定義 3.3 (1. 完全区別). 完全区別 PT_{dist} とは、 \mathcal{K} を言語族とするとき、 $PT_{dist}(\mathcal{K}, E) = 2^{\mathcal{K}} - \{\emptyset\}$ となる説明体系である。

完全区別の下でも、既知の言語すべて、つまり \mathcal{K} 自身を出力しなければならないとは限らないことに注意せよ。

3.3 単分化

次に挙げる説明体系は、外延の族が包含関係の列をなすような言語族のみを許すものである。ただし、言語同士が包含関係にあるとは限らない。直感的には、すべての例を説明する言語がありその一部分を次々と詳しく説明するという過程を定式化するものである。

定義 3.4 (2. 単分化). 説明体系 PT_{1-diff} が単分化であるとは、 $E \subseteq \Sigma^*$ とし、 \mathcal{K} を言語族とするとき、任意の $T \in PT_{1-diff}(\mathcal{K}, E)$ に対し、任意の $L, L' \in T$, $Ex(L, E) \subseteq Ex(L', E)$ または $Ex(L', E) \subseteq Ex(L, E)$. となる説明体系である。

3.4 分類

三番目の説明体系は、外延の族が交わりをもたないような言語族のみを許すものである。各言語どうしが共通部分を持たないとは限らない。直感的には、例をカバーしかつ分類するような言語族を許すものである。

定義 3.5 (3. 分類). 説明体系 PT_{group} が分類であるとは、 $E \subseteq \Sigma^*$ とし、 \mathcal{K} を言語族とするとき、任意の $T \in PT_{group}(\mathcal{K}, E)$, 任意の $L, L' \in T$ に対し、 $Ex(L, E) \cap Ex(L', E) = \emptyset$ となる説明体系である。

3.5 複分化

ここで挙げる最後の説明体系は、外延の族が包含関係の木 (正確には林) をなすような言語族のみを許すものである。直感的には、例の部分の説明する言語があれば、それを更に分類・説明する複数の言語の存在を許すようなものである。

定義 3.6 (4. 複分化). 説明体系 PT_{m-diff} が複分化であるとは、 $E \subseteq \Sigma^*$ とし、 \mathcal{K} を言語族とするとき、任意の $T \in PT_{m-diff}(\mathcal{K}, E)$ に対し、 $L, L' \in T$, $Ex(L, E) \cap Ex(L', E) \neq \emptyset$ ならば $Ex(L, E) \subseteq Ex(L', E)$ または $Ex(L', E) \subseteq Ex(L, E)$ となる説明体系である。

4 学習可能性に関する諸条件

4.1 完全区別に関する条件と例

このような出力をする場合には、以下が極限における区別の条件となる。

区別・対象で異なる言語族をとる場合

補題 4.1 (完全区別学習の必要十分条件). 言語 F が完全区別学習可能であるための必要十分条件は、 $\#\{F \cap L \mid L \in C\} < \infty$.

つまり、学習対象言語と区別言語の交差した結果が有限個になるということである。

任意の対象を学習する場合 次に, どのような言語族 \mathcal{C} により言語族 2^{Σ^*} が完全区別学習可能か考えてみよう. \mathcal{C} が無限言語族の場合には学習可能でない. 対象言語を $\Sigma^* \in 2^{\Sigma^*}$ にとると, 分割が無限に生じることが分かるからである. \mathcal{C} が有限であれば, 当然完全区別学習可能である.

完全区別学習可能な言語族は非常に限られている. 一例を挙げよう.

例 4.2. $L((m, n)) = \{(x, y) \in Q^2 \mid m \leq x \leq m+1, n \leq y \leq n+1, \}$, $\mathcal{L} = \{L((m, n)) \mid m, n \in N\}$ とし, \mathcal{F} を \mathcal{L} の言語の有限和からなる言語族, \mathcal{C} を \mathcal{L} の言語の和からなる言語族とすると, \mathcal{C} により言語族 \mathcal{L} は完全区別学習可能である.

否定的な条件 簡単のため, 対象となる言語族と区別する側の言語族が同一のものであるとき, どのような言語族が学習可能でないかについて議論しよう. 言語の正例からの学習において否定的な十分条件として知られている超有限な言語 (すべての有限言語と少なくとも一つの無限言語をもつような言語) は完全学習可能ではない. つまり正規言語の族に始まるチョムスキー階層の 4 つの言語族は完全区別学習可能でない.

正規言語より弱く, 有限言語より強い言語族については, 以下の条件を考慮するとよい.

条件: 語と言語の無限列で, $E_0 = \{e_0, e_1, \dots\}$, $E_n = \{e_n, e_{n+1}, \dots\}$, $\mathcal{L} = \{L_0, L_1, \dots\}$, $E_n \subseteq L_n$, $e_n \notin L_{n+1}$ を満たす.

対象となる言語が E_0 を含み, 区別する側の言語族が下記 \mathcal{L} を含む場合, 完全区別学習可能ではない. 対象となる言語族と区別する側の言語族が同一のものであるとき, 正規パターン言語族の非常に限られた部分族 $\{L(wx) \mid x \in X, w \in \Sigma^*\}$ ですら, 完全区別学習可能ではない.

4.2 単分化に関する条件

補題 4.3 (単分化の必要十分条件). F が \mathcal{C} で単分化学習可能であることの必要十分条件は, 任意の $\mathcal{L} \subseteq \mathcal{C}$ について, 言語 S_0, S_1, \dots, S_n が存在して, $\mathcal{E}_x(F, \mathcal{L}) = \{S_0, S_1, \dots, S_n\}$ かつ $S_0 \subseteq S_1 \subseteq \dots \subseteq S_n$ かつ任意の $L \in \mathcal{L}$ に対して $L \cap F \neq \emptyset$, ならば $\#\mathcal{L} < \infty$ となることである.

必要条件 言語族 \mathcal{C} により言語族 2^{Σ^*} が単分化学習可能であるための必要条件を二つ挙げておく. 第一の条件は有限の弾力性である.

定義 4.4 (有限の弾力性). [3] 言語族 \mathcal{C} が無限の弾力性をもつとは, 語の無限列 x_0, x_1, \dots と言語の無限列 $L_0, L_1, \dots \in \mathcal{C}$ が存在し, $\{x_0, x_1, \dots, x_n\} \subseteq L_n$ かつ $x_{n+1} \notin L_n$ ($n \in N$) となることである. また, \mathcal{C} が有限の弾力性をもつとは, \mathcal{C} が無限の弾力性をもたないことである.

第二の条件は, 下記条件の否定である.

言語族 \mathcal{C} が語の無限列 x_0, x_1, \dots と言語の無限列 $L_0, L_1, \dots \in \mathcal{C}$ が存在し, $\{x_n, x_{n+1}, \dots\} \subseteq L_n$ かつ $x_n \notin L_{n+1}$ ($n \in N$) となること.

有限の弾力性は言語の正例からの帰納推論の十分条件として知られるものである. 有限の弾力性をもたなくても, 正例からの帰納推論が可能な言語族も存在するため, 任意の言語を単分化できる言語族は, 正例からの帰納推論可能な族に真に含まれることになる.

4.3 分類に関する条件

補題 4.5 (分類の必要十分条件). F が \mathcal{C} で分類学習可能であることの必要十分条件は, 任意の $\mathcal{L} \subseteq \mathcal{C}$, 任意の $S, S' \in \mathcal{E}x(F, \mathcal{L})$ について, $S \cap S' = \emptyset$ かつ任意の $L \in \mathcal{L}$ に対して $L \cap F \neq \emptyset$, ならば $\#\mathcal{L} < \infty$ となることである.

必要条件 言語族 \mathcal{C} により言語族 2^{Σ^*} が分類学習可能であるための必要条件の一つとして, \mathcal{C} が言語の包含関係についての無限長の anti-chain を持たないことが挙げられる. 無限長 anti-chain とは以下のような性質である.

定義 4.6 (anti-chain). \leq が集合 A 上の半順序であるとする. \leq についての A 上の無限長 anti-chain とは, 無限列 $a_0, a_1, a_2, \dots, a_i, \dots$ で, 任意の a_i について, $i \neq j$ ならば a_i と a_j は \leq について比較不能であるようなものである.

4.4 複分化に関する条件

補題 4.7 (複分化の必要十分条件). F が \mathcal{C} で単分化学習可能であることの必要十分条件は, 任意の $\mathcal{L} \subseteq \mathcal{C}$ について, 以下の *i), ii)* が共に真となることである.

i) 言語 S_0, S_1, \dots, S_n が存在して, $\mathcal{E}x(F, \mathcal{L}) = \{S_0, S_1, \dots, S_n\}$ かつ $S_0 \subseteq S_1 \subseteq \dots \subseteq S_n$ かつ任意の $L \in \mathcal{L}$ に対して $L \cap F \neq \emptyset$, ならば $\#\mathcal{L} < \infty$ となることである.

ii) 任意の $S, S' \in \mathcal{E}x(F, \mathcal{L})$ について, $S \cap S' = \emptyset$ かつ任意の $L \in \mathcal{L}$ に対して $L \cap F \neq \emptyset$, ならば $\#\mathcal{L} < \infty$ となることである.

必要条件 言語族 \mathcal{C} により言語族 2^{Σ^*} が複分化学習可能であるためには, 少なくとも \mathcal{C} により 2^{Σ^*} が単分化学習可能であり, 分類学習可能でなければならない.

5 各説明体系下での学習可能性の関係

定理 5.1. F が言語, \mathcal{C} が帰納的言語の添え字つき族であるとき, *i) \Rightarrow ii), ii) \Rightarrow iii), iii) \Rightarrow iv), iii) かつ iv) \Rightarrow ii)* が成立する.

i) F が \mathcal{C} で完全区別学習可能である. *ii)* F が \mathcal{C} で複分化学習可能である. *iii)* F が \mathcal{C} で単分化学習可能である. *iv)* F が \mathcal{C} で分類学習可能である.

最後に本稿で挙げた学習モデルと他の学習モデルとの関係を挙げておく.

補題 5.2 (言語の非有界和の正例からの帰納推論の十分条件). [2] 言語族 \mathcal{L} が有限の弾力性を持ち, 言語の包含関係についての無限長 anti-chain をもたなければ, \mathcal{L} の非有界和からなる言語族は正例から帰納推論可能である.

定理 5.3. \mathcal{C} が帰納的言語の添え字つき族であるとき, *i)* 任意の言語が \mathcal{C} で単分化学習可能であるならば, \mathcal{C} は正例から帰納推論可能である. *ii)* 任意の言語が \mathcal{C} で分類学習可能であるならば, \mathcal{C} は言語の非有界和の正例からの帰納推論の十分条件を満たす.

6 結論

説明される入力例の分割の關係に注目した定式化により、各説明体系の關係、および正例からの単一言語・言語の非有界和の帰納推論との關係を明らかにした。

Reversible 言語族, Locally testable 言語族 等, 正規言語族より弱い既知の言語族がこの定式化でどのような位置に収まるかについての議論や, 最初に定義した 4 つの説明体系を拡張した. より応用向きと考えられる説明体系, 学習機械の出力はどのような性質をもち, 出力を利用する側にどのような恩恵をもたらすか等, 多くの議論すべき課題が残されているが, これらについては以降の研究課題としたい.

参考文献

- [1] E. M. Gold: *Language identification in the limit*, Information and Control, vol. 10, 447-474, (1967).
- [2] T. Shinohara and H. Arimura: *Inductive inference of unbounded unions of pattern languages from positive data*, Proc. the 7th International Workshop on Algorithmic Learning Theory, Lecture Notes in Artificial Intelligence, 1160, 256-271(1996).
- [3] K.Wright: *Identification of unions of languages drawn from positive data*, Proc. the 2nd Annual Workshop on Computational Learning Theory, 328-333, (1989)